Deflated Krylov Subspace Methods for
Nearly Singular Linear Systems[1]

by

Juan Camilo Meza and W.W. Symes[2]

Technical Report 87-3, February 1987.

| | | Form Approved<br>*OMB No. 0704-0188* |
|---|---|---|
| | **Report Documentation Page** | |

Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

| 1. REPORT DATE<br>**FEB 1987** | 2. REPORT TYPE | 3. DATES COVERED<br>**00-00-1987 to 00-00-1987** |
|---|---|---|
| 4. TITLE AND SUBTITLE<br>**Deflated Krylov Subspace Methods for Nearly Singular Linear Systems** | | 5a. CONTRACT NUMBER |
| | | 5b. GRANT NUMBER |
| | | 5c. PROGRAM ELEMENT NUMBER |
| 6. AUTHOR(S) | | 5d. PROJECT NUMBER |
| | | 5e. TASK NUMBER |
| | | 5f. WORK UNIT NUMBER |
| 7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)<br>**Computational and Applied Mathematics Department ,Rice University,6100 Main Street MS 134,Houston,TX,77005-1892** | | 8. PERFORMING ORGANIZATION REPORT NUMBER |
| 9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) | | 10. SPONSOR/MONITOR'S ACRONYM(S) |
| | | 11. SPONSOR/MONITOR'S REPORT NUMBER(S) |

12. DISTRIBUTION/AVAILABILITY STATEMENT
**Approved for public release; distribution unlimited**

13. SUPPLEMENTARY NOTES

14. ABSTRACT

15. SUBJECT TERMS

| 16. SECURITY CLASSIFICATION OF: | | | 17. LIMITATION OF ABSTRACT | 18. NUMBER OF PAGES | 19a. NAME OF RESPONSIBLE PERSON |
|---|---|---|---|---|---|
| a. REPORT<br>**unclassified** | b. ABSTRACT<br>**unclassified** | c. THIS PAGE<br>**unclassified** | | **20** | |

**Standard Form 298 (Rev. 8-98)**<br>Prescribed by ANSI Std Z39-18

# DEFLATED KRYLOV SUBSPACE METHODS FOR NEARLY SINGULAR LINEAR SYSTEMS *

J.C. MEZA AND W.W. SYMES †

**Abstract.** This paper concerns the use of Krylov subspace methods for the solution of nearly singular nonsymmetric linear systems. We show that the Incomplete Orthogonalization Methods (IOM) in conjuction with certain deflation techniques of Stewart and Chan can be used to solve large nonsymmetric linear systems which are nearly singular.

**1. Introduction.** This study concerns the use of Krylov subspace methods for the solution of nearly singular linear systems. There are many problems in numerical analysis which require the solution of nearly singular linear systems. For example, in the solution of nonlinear systems by homotopy continuation methods [1], or in nonlinear eigenvalue problems [5], one often has to solve nearly singular linear systems. Another example arises in constrained optimization problems [11], where the constraints may be nearly linearly dependent. Other examples include compartmental models [10], and decomposable Markov chains [16]. We focus on a problem from seismic processing known as the *velocity inversion problem* [18].

The goal of the velocity inversion problem is to determine certain parameters describing the earth from data taken at or near the surface. The data is known as a seismogram and we wish to determine the sound speed structure of the earth from these measurements. Given that we have a functional relation defined by $F(c) = s$, where $c$ is the speed of sound in the earth and $s$ is a seismogram, the inverse problem is: *given $s$ solve for $c$*. Since the seismogram is usually contaminated with noise we actually solve a nonlinear least squares problem — $\min J = \|s - F(c)\|$ for a suitable class of $c$.

A popular choice for the solution of nonlinear least squares problems is the Gauss-Newton method [8]. This method computes a sequence of iterates from the formulas

$$(1) \qquad (J(c_k)^T J(c_k))\Delta c = -J(c_k)^T F(c_k),$$

$$(2) \qquad c_{k+1} = c_k + \Delta c,$$

for $k = 0, 1, \ldots$, starting from an initial guess or model, $c_0$, for the sound speed structure. In the context of the velocity inversion problem the elements of the Jacobian matrix, $J(c_k)$, are not readily available. However, it can be shown that the Jacobian matrix acting on a vector can be computed from the solution of a certain boundary value problem. Similarly, the transpose of the Jacobian acting on a vector can be computed from the solution of a related boundary value problem. The solutions of the two boundary value problems are independent which results in a linear system

---

† Department of Mathematical Sciences, Rice University, P.O. Box 1892, Houston, Tx. 77251-1892.

which is almost symmetric. The Gauss-Newton Hessian, $J(c_k)^T J(c_k)$, would be symmetric except that the solution of the two boundary value problems, which define the Jacobian matrix and its transpose can only be solved up to some (nonsymmetric) discretization error.

The linear systems (1) that arise in the solution of the velocity inversion problem by the Gauss-Newton method are usually very large and since the matrix elements are not readily available the only recourse is to use an iterative technique for the solution of the linear systems. Additionaly, due to the physics of the problem the linear systems may be ill-conditioned either by having several large singular values or by having several small singular values. The small singular values arise because the initial data is band-limited. Perturbations in the model problem corresponding to frequencies outside the passband of the input data are simply not seen by the model. The large singular values are also inherent in this formulation of the problem. For more details the reader should consult [18].

In this study we address the issues of computing the solution to the linear systems (1) by using an iterative technique which computes the solution in the space spanned by the orthogonal complement of the singular vectors corresponding to the small singular values.

Several methods have been proposed for solving nearly singular linear systems. Consider the system of linear equations

$$(3) \qquad\qquad\qquad Ax = b,$$

where $x$ and $b$ are $n$ dimensional vectors and $A$ is an $n \times n$ real matrix which has rank $n - 1$. We denote the set of eigenvalues by $\lambda(A) = \{\lambda_1(A), ..., \lambda_n(A)\}$, with corresponding eigenvectors $\{u_1, u_2, ...u_n\}$. The eigenvalues are ordered so that

$$|\lambda_1| \geq |\lambda_2| \geq ... \geq |\lambda_n| = 0.$$

The solution to (3) can be written as

$$(4) \qquad\qquad x \;=\; x_d + \frac{u_n^T b}{\lambda_n}\, u_n,$$

$$(5) \qquad\qquad x_d \;=\; \sum_{i=1}^{n-1} \frac{u_i^T b}{\lambda_i}\, u_i,$$

where $u_n$ is a null vector of $A$. The vector $x_d$ is called the *deflated solution* to (3) and (4) is called the *deflated decomposition*. There are many definitions of the deflated solution. Chan [3] defines deflated solutions of (3) as solutions to nearby singular but consistent systems derived from (3). The choice of the nearby system will greatly affect the deflated solution. For example, one might choose the nearest singular matrix to $A$ in the Frobenius norm and pick the deflated solution to be the one with minimum norm. It is well known that this choice amounts to setting the smallest singular value of $A$ equal to zero in the singular value decomposition of the matrix $A$. Other examples can be found in [3].

2

In certain applications [4] it is preferable to compute the decomposition (4). In other applications, the deflated solution is the only solution of interest. Notice that if the eigenvector $u_n$ were known then both (4) and (5) could be computed by first computing $x$ and then orthogonalizing against $u_n$. However, even if the eigenvector $u_n$ were known this approach is not advisable because it usually results in a poor approximation to $x_d$ due to roundoff errors. In particular, if the component of the solution in the direction of the null eigenvector is large, then errors in that component tend to dominate the solution in the other directions.

Stewart [17] suggested a method for computing the deflated solution of (3) by an implicit method. This method uses orthogonal projections constructed from approximations to the singular vectors of the matrix $A$ corresponding to the smallest singular value. The disadvantage of this method is that it requires a direct method for the solution of (3). Chan [6] proposed a deflated Lanczos method for symmetric positive definite linear systems which only requires the matrix-vector product $Ax$. The goal of this paper is to study methods for computing the deflated solution for large sparse nonsymmetric linear systems which are nearly singular.

The basic idea is to use one of the Krylov subspace methods proposed by Saad [14] for solving nonsymmetric linear systems. These methods compute an approximation to the solution by generating iterates which lie in a certain Krylov subspace. The approximations to the solution are then computed by solving a linear system described by a small upper Hesssenberg matrix, $H$, produced by Arnoldi's method [2].

Arnoldi's method may be thought of as a Galerkin process for approximating the eigenvalues of $A$ by the eigenvalues of the upper Hessenberg matrix $H$. Like the Lanczos method, the approximations to the eigenvalues tend to be best at the extremes of the spectrum, so that the matrix $H$ may be expected to also be nearly singular. Therefore, one disadvantage to using the methods proposed by Saad is that they could require the solution of a nearly singular linear system. Unlike the original system (3) however, we need only be able to solve linear systems involving the much smaller matrix $H$. We can then use the deflation techniques suggested by Stewart and Chan. In addition we propose another method based on solving a truncated least squares problem analogous to the GMRES method suggested by Saad [15].

**2. Krylov Subspace Methods.** Saad [14] proposed a class of methods for solving large sparse nonsymmetric linear systems based on the Arnoldi process [2] for computing the eigenvalues of a matrix. Arnoldi's method is a generalization of the Lanczos method [13] for nonsymmetric matrices, and when it is applied to a symmetric matrix it reduces to the Lanczos method. Like the Lanczos method, Arnoldi's method is best viewed as an iterative method for approximating the eigenvalues of large sparse matrices. In essence, Arnoldi's method is just the Gram-Schmidt method for computing an orthonormal basis for the Krylov subspace $\kappa_m(w_1, A) \equiv \text{span}\{w_1, Aw_1, ..., A^{m-1}w_1\}$. The method can by described as follows.

ALGORITHM 2.1. *Arnoldi's Method.*
1. Choose $w_1$ such that $\|w_1\| = 1$.
2. For $j = 1, 2, \dots$

$$
\begin{aligned}
h_{ij} &= (Aw_j, w_i) \qquad i = 1, 2, \dots, j \\
\hat{w}_{j+1} &= Aw_j - \sum_{i=1}^{j} h_{ij} w_i \\
h_{j+1,j} &= \|\hat{w}_{j+1}\| \\
w_{j+1} &= \hat{w}_{j+1} / h_{j+1,j}
\end{aligned}
$$

If we let $W_m = [w_1, w_2, \dots, w_m]$, then it is easy to show that

$$(6) \qquad H_m = W_m^T A W_m,$$

where the entries of the upper Hessenberg matrix $H_m$ are the scalars $h_{ij}$ produced by Arnoldi's method after $m$ steps.

As Saad [14] has shown, Arnoldi's method may be used as a basis for a class of Krylov subspace methods for solving large sparse nonsymmetric linear systems. By a Krylov susbpace method, we mean any method that approximates the solution to the linear system (3) by generating iterates of the form

$$x_m = x_0 + z_m,$$

where $x_0$ is an initial guess, $z_m \in \kappa_m(r_0, A)$ and $r_0 = b - Ax_0$. If we carry out $m$ steps of Arnoldi's method starting with $w_1 = r_0/\|r_0\|$ and if we impose the Galerkin condition that the residuals at each iteration be orthogonal to $\kappa_m(r_0, A)$, then this yields

$$W_m^T A W_m y_m - W_m^T r_0 = 0.$$

Using the relation (6) yields

$$x_m = x_0 + W_m y_m,$$

where $y_m$ solves the system

$$(7) \qquad H_m y_m = \|r_0\| \, e_1,$$

and $e_1 = (1, 0, \dots, 0)^T$.

This defines the Full Orthogonalization Method [14], for solving (3).

ALGORITHM 2.2. *Full Orthogonalization Method.*
  1. Choose $x_0$ and compute $r_0 = b - Ax_0$. Set $w_1 = r_0 / \|r_0\|$.
  2. For $j = 1, 2, ...m$

$$
\begin{aligned}
h_{ij} &= (Aw_j, w_i) \qquad i = 1, 2, ..., j \\
\hat{w}_{j+1} &= Aw_j - \sum_{i=1}^{j} h_{ij} w_i \\
h_{j+1,j} &= \|\hat{w}_{j+1}\| \\
w_{j+1} &= \hat{w}_{j+1} / h_{j+1,j}
\end{aligned}
$$

  3. Form the solution:
    - Solve $H_m y_m = \|r_0\| e_1$
    - Set $x_m = x_0 + W_m y_m$.

In practice, the number of iterations is chosen so that the approximate solution $x_m$ is sufficiently accurate. Usually this is measured by requiring that the initial residual be decreased by a user-specified amount. Fortunately, the residual at any iteration may be computed without actually computing the solution to (3) through the relation [14],

$$
(8) \qquad \|b - Ax_m\| = h_{m+1,m} |e_m^T y_m|.
$$

Although the computation of the residual by (8) requires solving (7) for $y_m$ there are ways to circumvent this computation by carrying an LU or QR factorization of $H$ throughout the Arnoldi process. If after $m$ iterations the approximate solution has not converged then it is possible to restart the algorithm using the current estimate of $x$ as the new initial guess. This method is denoted by FOM($k$), or the restarted FOM.

It is well known that the Arnoldi process may be viewed as a Galerkin process for estimating the eigenvalues of a matrix. In particular, if we apply Algorithm 2.2 to a linear system that is nearly singular then the upper Hessenberg matrix, $H_m$, which is generated after m steps of the Arnoldi process will probably have a small eigenvalue. Therefore if we solve (7) for $y_m$ in the straightforward way, our computed solution will be inaccurate for the reasons indicated in Section 1.

Fortunately, computing the deflated solution of (7) is easier than computing the deflated solution of (3). Since the matrix $H_m$ has dimension $m \ll n$ the solution of (7) is at least computationally easier. Moreover, the matrix elements of $H_m$ are on hand whereas the matrix elements of $A$ are not available. The next section describes some of the deflation techniques which can be used to compute the deflated solution of (7) in a stable manner.

**3. Deflation Methods.** This section describes three methods for computing the deflated solution to a nearly singular linear system. In particular, we consider the

linear system

$$(9) \qquad\qquad H_m y = f,$$

where $y$ and $f$ are $m$ dimensional vectors and $H_m$ is the upper Hessenberg matrix generated after $m$ steps of Arnoldi's method. Further, we assume that $H_m$ is nearly singular with a rank deficiency of at most one. The extensions to null spaces of dimensions greater than one will be treated in a later section.

The first method, which uses orthogonal projections, was proposed by Stewart [17]. The second method is a generalization of a technique suggested by Chan [6] for symmetric positive definite systems. In essence, this method uses a QR iteration to decouple the linear system into a well conditioned problem plus a component corresponding to the small eigenvalue. The third method computes the deflated solution by solving a truncated least squares problem.

**3.1. Deflation by Orthogonal Projection.** Stewart [17] proposed an implicit method for computing the deflated solution to a nearly singular linear system. His algorithm consists of constructing two orthogonal projectors defined by approximations to the singular vectors corresponding to the small singular value.

Consider the right and left singular vectors, respectively $v$ and $u$, corresponding to the smallest singular value of $H_m$. Define the orthogonal projectors $P_u = I - uu^T$ and $P_v = I - vv^T$. The projector $P_u(P_v)$ is merely the orthogonal projector onto the orthogonal complement of the space spanned by $u(v)$. It is easy to show that the deflated solution to (9) is the unique vector satisfying the relations

$$(10) \qquad\qquad P_u H_m P_v\, y_d \;=\; P_u f,$$
$$(11) \qquad\qquad P_v\, y_d \;=\; y_d.$$

Stewart suggests using the following algorithm based on iterative refinement to solve for $y_d$ .

> ALGORITHM 3.1. *Deflation by Orthogonal Projection.*
> 1. Set $y = 0$
> 2. For $k = 1, 2, \ldots$
>    - Solve $H_m d = P_u(f - H_m y)$,
>    - Set $y = y + P_v d$.
> 3. $y_d = y$

Since the singular vectors are not known, Stewart suggests approximating $u$ and $v$ by a variant of the inverse power method. In the case of interest, where the small singular value is isolated the inverse power method is known to converge rapidly. Stewart also gives conditions, which depend on the accuracy attained in the approximation to the vectors $u$ and $v$, under which Algorithm 3.1 will converge to the deflated solution of (9).

**3.2. Deflation by QR iteration.** The second method for computing the deflated solution is a generalization of a method proposed by Chan [6] for symmetric positive definite systems. To motivate the discussion, assume that we have the pair $(\lambda_m, u_m)$ of the unreduced upper Hessenberg matrix $H_m$. It is well known that one step of the shifted QR iteration method with a shift of $\lambda_m$ reduces the matrix $H_m$ so that the eigenvalue $\lambda_m$ appears on the diagonal and the corresponding subdiagonal element is zero. The linear system (9) can then be decoupled so that it is easily solved for the deflated solution.

In particular, if we compute the matrix $H^{(1)}$ by the following two step procedure: Compute the QR factorization of $H_m - \lambda_m I$ Form $H^{(1)} = RQ + \lambda_m I$, then the matrix $H^{(1)}$ takes on the form

$$H^{(1)} = \begin{array}{c} m-1 \\ 1 \end{array} \begin{pmatrix} \overset{m-1}{\hat{H}} & \overset{1}{\hat{h}} \\ 0 & \lambda_m \end{pmatrix},$$

where $\hat{H}$ is an upper Hessenberg matrix of order $m - 1$. It is easy to show that $H^{(1)} = Q^T H_m Q$ , so that we can transform (9) into

$$(12) \qquad\qquad H^{(1)}z = \hat{f},$$

where $y = Qz$ and $\hat{f} = Q^T f$.

If we partition the vectors $z$ and $\hat{f}$ so that

$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}, \qquad \hat{f} = \begin{bmatrix} \hat{f}_1 \\ \hat{f}_2 \end{bmatrix},$$

then the deflated solution to (12) is

$$(13) \qquad\qquad z_d = \begin{bmatrix} \hat{H}^{-1}\hat{f}_1 \\ 0 \end{bmatrix}.$$

The deflated solution to (9) can now be computed from $y_d = Qz_d$.

There are two remarks in order. The first is that in general we do not know the eigenvalue $\lambda_m$; however, as in the previous section, we may compute an approximation to $\lambda_m$ by using the inverse power method. The second remark is that even if we had an exact value for, $\lambda_m$, in practice the matrix $H_m$ will not be reduced in one QR iteration step. Wilkinson [19] suggests iterating with the shifted QR method until the element on the subdiagonal has converged to zero. We use the test

$$(14) \qquad\qquad \|H^{(1)}_{m,m-1}/H^{(1)}_{m,m}\| \le \epsilon$$

to check for convergence. Two or three iterations are usually sufficient to satisfy (14). This leads to the following algorithm.

ALGORITHM 3.2. *Deflation by QR iteration.*
1. Compute $\lambda_m$ from $H_m$ via inverse iteration.
2. $H^{(0)} = H_m$
3. For $j = 0, 1, \ldots$ until convergence.
    - Compute the QR factorization of $H^{(j)} - \lambda_m I$.
    - Compute $H^{(j+1)} = RQ + \lambda_m I$.
4. Set $z_d = \left[ \hat{H}^{-1} \hat{f}_1, 0 \right]^T$.
5. Set $y_d = Q z_d$.


If the dimension of the null space is greater than one then the issues become more complicated. There are two cases to consider: a small eigenvalue of multiplicity greater than one and the case of several distinct small eigenvalues. The case of a small real eigenvalue with a multiplicity greater than one is easily handled since the shifted QR method will still converge. The case of several distinct small eigenvalues is harder to address since the QR method does not guarantee that the converged eigenvalues will appear in any particular order on the diagonal of $H^{(1)}$. A related issue is that of complex eigenvalues. Since the inverse power method converges to the eigenvalue of smallest modulus, we must be careful in choosing the shift. A good choice would be to use Francis' implicit double shift QR method and with a shift as described by Wilkinson [19].

**3.3. Deflation by Truncated Least Squares.** The last method we discuss is derived from a technique for solving rank deficient least squares problems [12, pp. 162-167]. The general idea in these problems is to compute a QR factorization of the nearly singular matrix $H_m$ such that the elements on the diagonal of the matrix $R$ display the rank deficiency. The solution to the linear system (9) is then computed by setting the small diagonal elements of the matrix $R$ equal to zero and solving a *truncated least squares* problem.

Assume that the matrix $H_m$ has exactly one zero eigenvalue and consider its QR factorization,

$$H_m \Pi = QR,$$

where $\Pi$ is a permutation matrix chosen so that the matrix $R$ has the form

$$R = \begin{array}{c} m-1 \\ 1 \end{array} \begin{array}{cc} m-1 & 1 \\ \left( \begin{array}{cc} R_{11} & R_{12} \\ 0 & 0 \end{array} \right), \end{array}$$

where $R_{11}$ is upper triangular.

In this case, the least squares solution to (9) can be easily computed. In fact,

$$\| H_m y - f \|^2 = \| R_{11} z_1 - (\hat{f}_1 - R_{12} z_2) \|^2 + \| \hat{f}_2 \|^2,$$

where

$$(15) \qquad \qquad \Pi^T y = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

and

$$(16) \qquad Q^T f = \begin{bmatrix} \hat{f_1} \\ \hat{f_2} \end{bmatrix}.$$

If we set $z_2 = 0$, then the solution is given by

$$(17) \qquad \hat{y} = \Pi \begin{bmatrix} R_{11}^{-1} \hat{f_1} \\ 0 \end{bmatrix}.$$

The vector $\hat{y}$ is called the basic solution. In general, it is not the least norm least squares solution unless $R_{12} = 0$.

In practice, the matrix $H_m$ will never have an eigenvalue exactly equal to zero. The question is then to determine the rank of the matrix $H_m$ from the elements of $R$. One popular choice is the method of column pivoting implemented in LINPACK [9]. This method is usually reliable in detecting the rank deficiency of a matrix, although there are counterexamples where the method may fail. Chan [7] suggested another method for producing a $QR$ factorization which guarantees a small $R_{22}$ element. The essential ideas can be explained in the case of a rank one deficient matrix. First, we need the following theorem proved by Chan [7].

THEOREM 1. *Suppose that $x \in \mathbb{R}^n$ with $\|x\| = 1$ such that $\|Ax\| = \epsilon$. Let $\Pi$ be a permutation such that if $\Pi^T x = y$, then $|y_n| = \|y\|_\infty$ . Then if $A\Pi = QR$ is the $QR$ factorization of $A\Pi$, then*

$$|r_{nn}| \leq \sqrt{n}\,\epsilon.$$

The usefulness of this theorem is apparent if we consider the right singular vector, $v$, of the matrix $H_m$ corresponding to the smallest singular value $\sigma_m$. Then we have $\|v\| = 1$, and $\|H_m v\| = \sigma_m$. If we define the permutation $\Pi$ by

$$|(\Pi^T v)_m| = \|v\|_\infty,$$

then $H_m \Pi = QR$ has a pivot $r_{mm}$ at least as small as $\sqrt{m}\,\sigma_m$ in absolute value. As in Stewart's method, all that is required is an approximation to $v$, which may be computed by the inverse power method. This suggests the following algorithm for solving (9).

ALGORITHM 3.3. *Deflation by Truncated Least Squares.*
1. Compute the QR factorization of $H_m$.
2. Compute $v$ and $\sigma_m$ from $H_m$ via inverse power method.
3. Compute $\Pi$ so that $|(\Pi^T v)_m| = \|v\|_\infty$
4. Compute the QR factorization of $H_m \Pi$.
5. Compute $\hat{y} = \Pi \left[ R_{11}^{-1} \hat{f_1}, 0 \right]^T$.

9

This method has the obvious advantage of being immediately applicable to null spaces with a dimension greater than one. However the question of rank deficiency is still a hard problem and the user must be able to supply a tolerance which specifies the amount of ill-conditioning allowed.

**4. Numerical Results..** This section presents several numerical experiments comparing the various methods described in Section 3.

Recall that he linear system of interest is

$$Ax = b,$$

where $x$ and $b$ are $n$ dimensional vectors and $A$ is an $n \times n$ real matrix which is nearly singular. All of the numerical results presented here are for linear systems of order 10. The results are similar for larger systems. The numerical experiments were run on a Sun 3/160 computer, using single precision arithmetic (machine epsilon $\approx 10^{-7}$). The method was said to converge whenever

$$\frac{\|r_k\|}{\|r_0\|} \leq 10^{-6}.$$

The basic method used was the restarted $FOM(k)$. To avoid the issue of deciding when the eigenvalues of the upper Hessenberg matrix, $H$, are good approximations to the eigenvalues of the matrix A, we set $k = n$.

The first test case was a small peturbation to a symmetric positive definite matrix. Define

$$A(\epsilon) = D + \epsilon E.$$

The matrix $D$ is defined by $D = diag\ (10^{-I}, 2, 3, \ldots, n)$, and $I$ varies from 1 to 7 . The nonsymmetric perturbations are generated using the random number generator, URAND, from IMSL. The matrices, $E$, are computed by generating uniform random numbers between [-0.5, +0.5 ], and normalizing so that $\|E\|_2 = 1$. The amount of nonsymmetry can then be adjusted by varying $\epsilon$. We set $\epsilon = 10^{-3}$.

The second set of test cases was picked from a study done by Chan [6]. In this set of problems we set

$$A = (I - uu^T)\, D\, (I - vv^T),$$

where $u$ and $v$ were chosen randomly with the constraint that they have norm one. The matrix D was chosen as in problem 1.

The purpose of the third problem is to simulate a typical linear system arising in the velocity inversion problem. These problems usually have one or more small singular values and one or more large singular values with the rest of the spectrum fairly well-conditioned. In this problem, we set the matrix $D = diag\ (10^{-I}, 1, \ldots, 3, 3000)$ , with the values of $d_2$ through $d_{n-1}$ varying uniformly between 1 and 3. The matrix A is then computed as in problem 1. This example generates a well-conditioned problem if the small and large eigenvalues are excluded, which is typical of some of the velocity inversion problems.
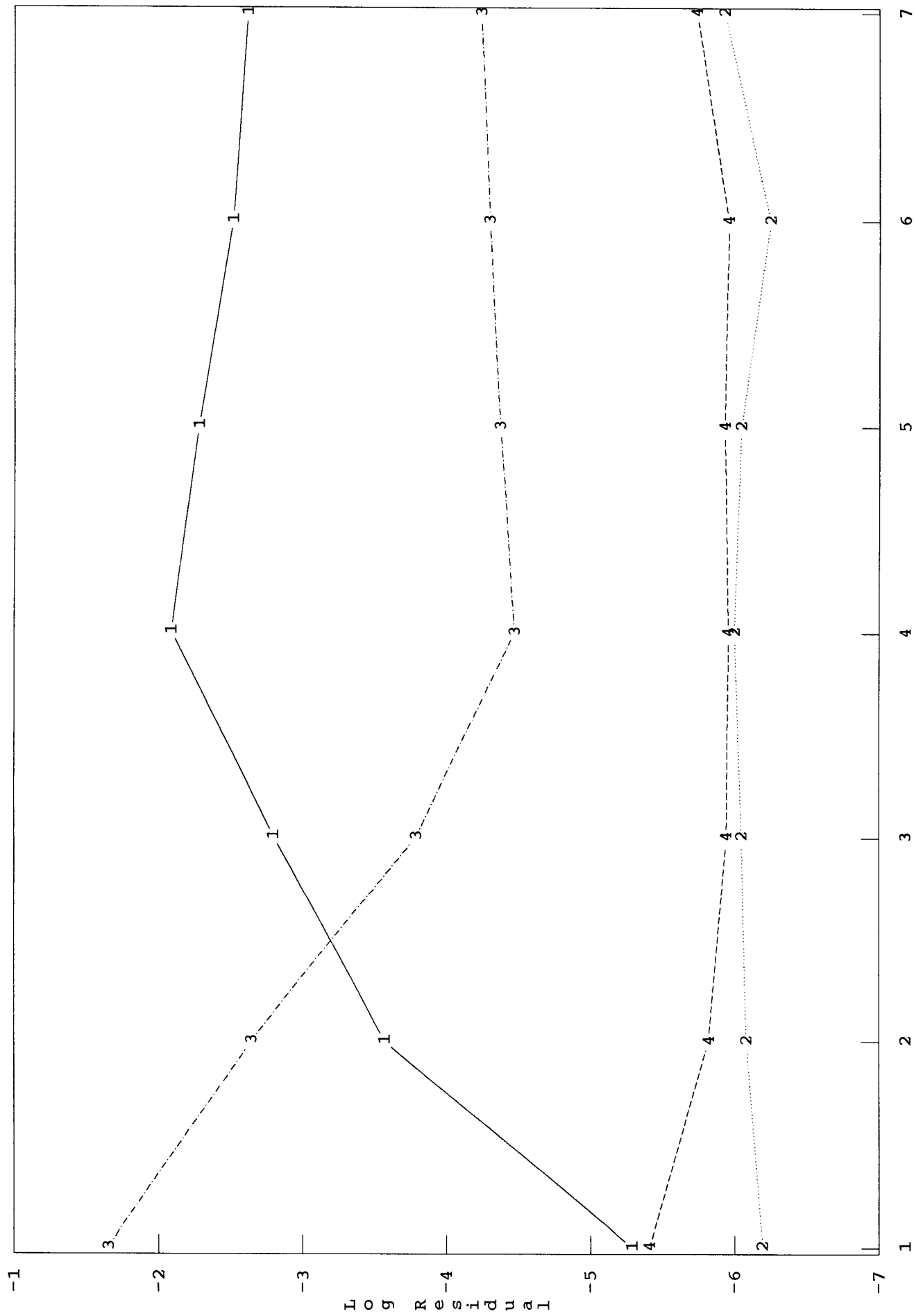
The four methods discussed in Section 3 are:

10

1. Full Orthogonalization Method with deflation after the solution of Ax = b. — FOM.
2. FOM with Stewart's orthogonal projection deflation — FOMOP.
3. FOM with the truncated least squares deflation — FOMLSQ.
4. FOM with the QR iteration deflation — FOMQRI.

Figures 1-3 plot the $\log\|r\|$ versus $I$. Since we are interested in the deflated solution we have chosen to plot the norm of the residual corresponding to the deflated solution. As the value of $I$ increases, the linear systems become more ill-conditioned, with the smallest singular value of each matrix approximately equal to $10^{-I}$. The effect of the near singularity of the linear systems on Method 1 is apparent as the plots show the norm of the residual increasing as the linear system becomes more ill-conditioned. This is to be expected as the error in the component corresponding to the smallest singular values tends to contaminate the rest of the solution. Method 2, FOMOP, is clearly the best method in terms of producing the smallest residuals. Method 3, FOMLSQ, was disappointing in that the norm of the residual was consistently worse than the other methods. Method 4, FOMQRI, was inconsistent in this set of problems. In problems 1 and 3, FOMQRI was almost as good as FOMOP. However in problem 2 the method performed much worse.

The question of extending these methods to null spaces of higher dimensions is also of interest. In this respect, FOMLSQ can be easily extended while the other methods would require some extra work. In fact the ease with which FOMLSQ can be extended to solving systems which have several small singular values is perhaps the only redeeming factor of this method. Figures 4-6 illustrate the effect of solving the same systems of problems 1-3 using both a truncated least squares approach and simply solving the system using a full least square problem. It can be seen that the two curves usually intersect between $I = 2$ and $I = 3$ which corresponds to $\kappa(A) \approx 1000$. This implies that using FOMLSQ with a properly chosen tolerance would perform better than always solving the truncated least squares problem. Fortunately this is the case of interest. This approach would still not be as good as using FOMOP, but it has the advantage that it is easier to extend to null spaces with dimension greater than one, whereas extending FOMOP to handle several small singular values would require extra work.

Further experiments remain to be done. It is not clear which deflated solution is best in terms of the norm of the error. In fact, there are many definitions of the deflated solution and the choice of the deflated solution greatly alters the results. That is one reason we have chosen to work with the norm of the residual. Another possibility lies in solving the upper Hessenberg system by using the singular value decomposition. Although this may seem like to much work at first sight, the problems we are dealing with have the property that one matrix-vector multiply is very expensive. In this context, computing the SVD for a small upper Hessenberg matrix would be insignificant. This approach also has the added attraction that (like the FOMLSQ method) it is easily extended to null spaces of dimension greater than one. This research will be the subject of a later report.

Figure 1. Problem 1 --- Log || r || Vs. I

1=FOM, 2=FOMOP, 3=FOMLSQ, 4=FOMQRI
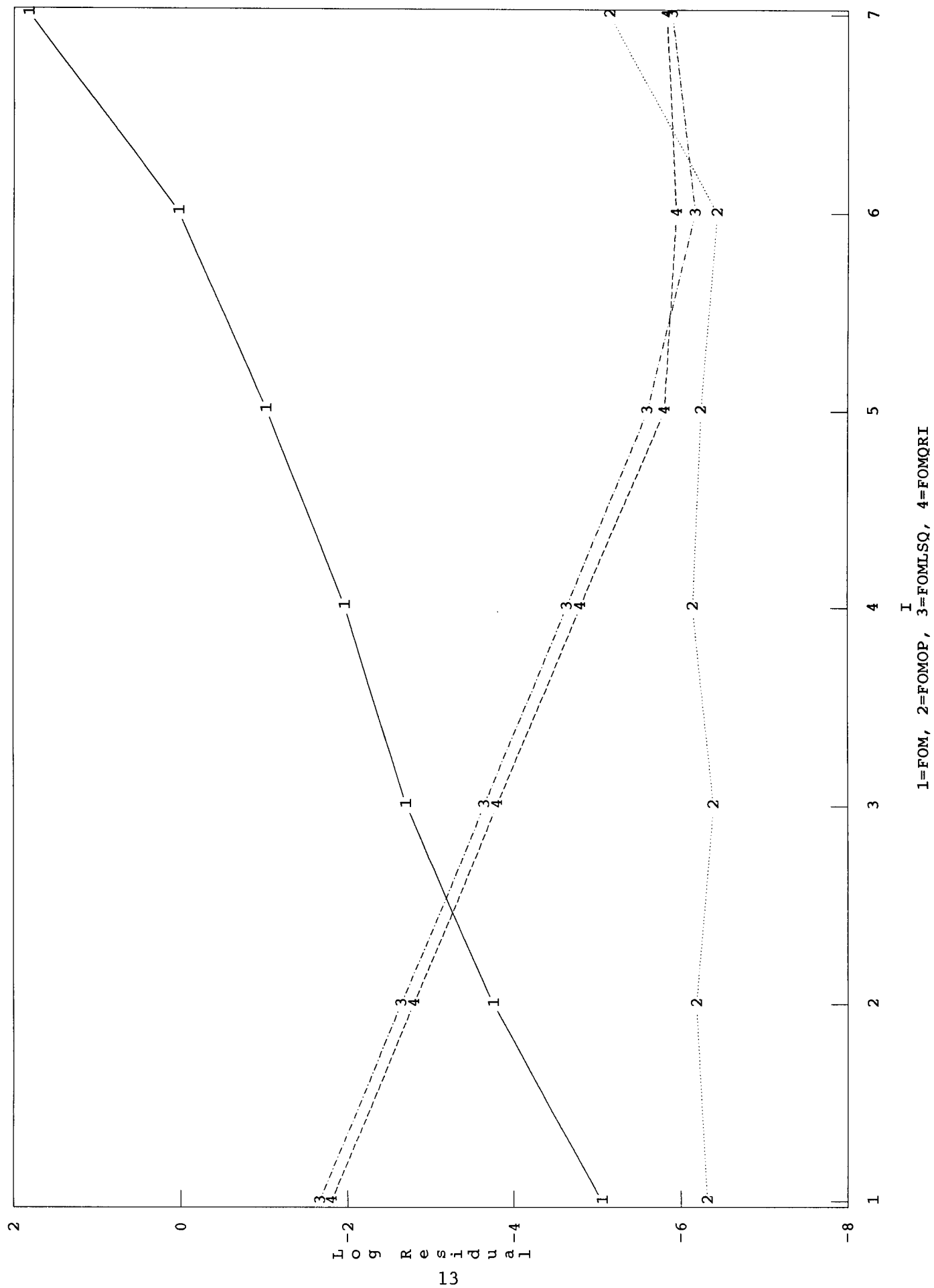
12

Figure 2. Problem 2 --- Log || r || Vs. I

1=FOM, 2=FOMOP, 3=FOMLSQ, 4=FOMQRI

13

Figure 3. Problem 3 --- Log || r || Vs. I

1=FOM, 2=FOMOP, 3=FOMLSQ, 4=FOMQRI

14

Figure 4. Problem 1 --- Effect of Truncation on FOMLSQ

I
1=FOMLSQ, 2=FOMLSQ w/o Truncation

15

Figure 5. Problem 2 --- Effect of Truncation on FOMLSQ

1=FOMLSQ, 2=FOMLSQ w/o Truncation

16

Figure 6. Problem 3 --- Effect of Truncation on FOMLSQ

1=FOMLSQ, 2=FOMLSQ w/o Truncation

17

## REFERENCES

[1] E. ALLGOWER AND K. GEORG, *Simplicial and continuation methods for approximating fixed points and solutions to systems of equations*, SIAM Rev., 22 (1980), pp. 28–85.

[2] W. ARNOLDI, *The principle of minimized iteration in the solution of the matrix eigenvalue problem*, Quart. Appl. Math., 9 (1951), pp. 17–29.

[3] T. CHAN, *Deflated decomposition of solutions of nearly singular systems*, SIAM J. Numer. Anal., 21 (1984), pp. 738–754.

[4] ———, *Deflation techniques and block-elimination algorithms for solving bordered singular systems*, SIAM J. Sci. Stat. Comput., 5 (1984), pp. 121–134.

[5] ———, *Newton-like pseudo-arclength methods for computing simple turning points*, SIAM J. Sci. Stat. Comput., 5 (1984), pp. 135–148.

[6] ———, *Deflated Lanczos procedures for solving nearly singular systems*, Research Report YALEU/DCS/RR-403, Department of Computer Science, Yale University, 1985.

[7] ———, *Rank revealing QR factorizations*, Research Report YALEU/DCS/RR-398, Department of Computer Science, Yale University, 1985.

[8] J. DENNIS AND R. B. SCHNABEL, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*, Prentice-Hall Inc., Englewood Cliffs, New Jersey, 1983.

[9] J. DONGARRA, J. BUNCH, C. MOLER, AND G. STEWART, *LINPACK Users' Guide*, SIAM Press, Philadelphia, PA, 1984.

[10] R. FUNDERLIC AND J. MANKIN, *Solution of homogeneous systems of linear equations arising from compartmental models*, SIAM J. Sci. Stat. Comput., 2 (1981), pp. 375–383.

[11] P. GILL, W. MURRAY, AND M. WRIGHT, *Practical Optimization*, Academic Press, New York, 1979.

[12] G. GOLUB AND C. VAN LOAN, *Matrix Computations*, Johns Hopkins University Press, Baltimore, MD, 1983.

[13] C. LANCZOS, *An iteration method for the solution of the eigenvalue problem of linear differential and integral operators*, J. Res. Nat. Bureau Standards, 45 (1950), pp. 255–282.

[14] Y. SAAD, *Krylov subspace methods for solving large unsymmetric linear systems*, Math. Comp., 37 (1981), pp. 105–126.

[15] Y. SAAD AND M. SCHULTZ, *GMRES: a generalized minimal residual algorithm for solving nonsymmetric linear systems*, SIAM J. Sci. Stat. Comput., 7 (1986), pp. 856–869.

[16] G. STEWART, *Computable error bounds for aggregated Markov chains*, Technical Report 901, University of Maryland Computer Science Center, College Park, MD, 1980.

[17] G. STEWART, *On the implicit deflation of nearly singular systems of linear equations*, SIAM J. Numer. Anal., 2 (1981), pp. 136–140.

[18] W. SYMES, *Stability properties for the velocity inversion problem*, in SEG/SIAM/SPE Symposium, Houston, TX, January 1985.

[19] J. WILKINSON, *The Algebraic Eigenvalue Problem*, Oxford University Press, London, 1965.